# STATIC END TO END RETRANSMIT APPARATUS AND METHOD

1 **Technical Field**

2       The technical field is error detection and correction in multiprocessor computer

3 systems.

4 **Background**

5       Path or link errors may exist in multiprocessor computer systems. To tolerate

6 such link errors, computer designers have traditionally made use of error correction code

7 (ECC) or retry mechanisms. ECC handles certain permanent errors such as a wire being

8 disconnected in a link (or interconnect) while other links are working. However, if

9 multiple wires in the link are disconnected, or if the entire link is disconnected, the ECC

10 cannot recover the disconnected link. Retry works well for transient errors. If a packet

11 includes errors that can be detected, but not corrected, then the packet will be sent again

12 from a sending node to a receiving node using the same link. The process of sending the

13 packet may repeat several times. However, retry cannot handle errors such as multiple

14 wires failing in a link or the link being disconnected, or an intermediate routing chip

15 being removed for service.

16       An end to end retry scheme may be used as a means to tolerate link or immediate

17 route chip failures. The basic approach is that each transaction has a time to live, and as a

18 transaction travels through the multiprocessor computer architecture, the value of the time

19 to live is decremented. A transaction that cannot be delivered to its destination node and

20 has its time to live go from its maximum level to zero is discarded. Request transactions

21 may be retried along a secondary path if some predetermined number of attempts along

22 the primary path failed to generate a response. Response transactions may not be

23 acknowledged. If a response transaction does not reach its destination mode, the failure

24 of the response transaction to reach the destination node will have the same effect as a

25 request transaction not reaching the destination mode, and as a result the request

26 transaction may be retried.

27       This end-to-end retry scheme has several disadvantages. First, is that the time-out

28 hierarchy is tied to the retry protocol. If a request transaction is tried four times, for

29 example, (along primary and alternate paths) before the request reaches an error time out,

30 then the next transaction type in the hierarchy has to wait for four times the time out for

31 every lower level transaction, the transaction type can generate. For example, a memory

32 read request may cause several recalls. Thus, the memory read request may be reissued

1   only after allowing all recalls to happen. Thus, the memory read request's reissue time

2   out is the maximum number of recalls times the four times the recall time out, plus the

3   times of flight for the request transaction and the response transaction. As a result, the

4   time out hierarchy keeps increasing exponentially (that is the factor four keeps getting

5   multiplied across the hierarchy).

6       A second disadvantage is that verifying a time out hierarchy is a challenging

7   design requirement since time outs frequently take place over the period of time measured

8   in seconds, and simulating a large system to the range of seconds of operation is almost

9   impossible. A third disadvantage is that the retry strategy requires participation of all

10  chips in the interconnect (at least to decrement the time out value). Thus, the retry

11  strategy does not work well in a computer architecture that has components, such as a

12  crossbar, that the computer designer is trying to leverage that is oblivious to the recovery

13  features such as time to live. A fourth disadvantage is that the retry strategy operates in

14  an unordered network, and ordered transactions such as processor input/outputs (PIOs)

15  need an explicit series of sequence numbers to guarantee ordering. In addition, for

16  transactions such as PIO reads that have side effects, a read return cache is needed to

17  ensure the same PIO read is not forwarded to a PCI bus multiple times.

18  **Summary**

19      A static end to end retry apparatus and method uses the concept of sequence

20  numbers combined with a path number. All transactions sent along a path are delivered

21  in order to remove any time dependency. The apparatus and method ensure there are no

22  duplicate transactions through the use of special probe and plunge transactions and their

23  respective responses. The apparatus and method also allow for any number of alternate

24  paths being active simultaneously, such that if one path fails, the remaining alternate

25  paths can continue on the communication (along with the backup alternate path if desired)

26  as usual without any loss of transactions.

27      In a multiprocessor computer system with multiple nodes, each node keeps track

28  of transactions the node has sent over time to every other node, as well as every

29  transaction the node has received from every other node along each active path for each

30  flow control class. To accomplish this tracking function, two data structures exist. A

31  send_seqid, representing the sequence identification (ID) (or sequence number) for the

32  last transaction sent by the sending (or source) node to a given destination node exists

33  along any given active path, and a flow control class. A second structure is a

34  receive_seqid, representing the sequence ID (sequence number) of the last transaction that

1     a destination node received and for which the destination node sent an acknowledgement

2     (ACK) back to the source node, for each node, along every active path, and for each flow

3     control class. The send_seqid and the receive_seqid may be stored in send_seqid and

4     receive_seqid tables at each node in the multiprocessor computer system.

5     Each node (destination node for the send_sequid or source node for the

6     receive_sequid) can receive transactions over multiple paths. All nodes in one flow

7     control class may operate over the same number of paths. For example, the system may

8     have four alternate active paths between any two CPU/memory nodes, but only one active

9     path to or from an I/O hub chip. The system does not require distinct physical paths

10     between any source-destination nodes. For example, the system may comprise four

11     active paths with two active paths sharing a physical path.

12     Every transaction that is sent from a source node to a destination node is also put

13     into a retransmit buffer. When the transaction results in an acknowledgement from the

14     destination node, the transaction is removed from the retransmit buffer. The

15     acknowledgement can be piggy-backed with an incoming transaction and/or a special

16     transaction. No acknowledgement is necessary for an acknowledgement transaction. If a

17     transaction is not acknowledged within a predetermined time, recovery actions are taken.

18     The destination node may wait to gather several transactions for a given source node

19     before generating an explicit acknowledgement transaction, while trying to ensure that

20     such a delay will not generate any recovery actions at the source node. This delay helps

21     conserve bandwidth by avoiding explicit acknowledgement transactions as much as

22     possible.

23     When a source node sends a transaction to a destination node, the source node

24     gets the sequence number from the send_seqid table, checks that no transaction with the

25     source sequence number is pending to the same destination node, and sends the

26     transaction to the destination node while also sending the transaction to the retransmit

27     buffer. The source node then increments the corresponding sequence number in the

28     send_seqid table. When the destination node receives the transaction, the destination

29     node queues the transaction in a destination node buffer. If the transaction is of a request

30     type, and the destination node can generate a response within a time out period, the

31     destination node sends a response, which acts as in implicit acknowledgement, to the

32     source node. The destination node checks the receive_seqid table to see if the transaction

33     received by the destination node has the correct sequence number. If the transaction has

34     the correct sequence number, the destination node updates the corresponding entry in the

1　receive_seqid table, and sets a flag bit indicating that the destination node needs to send

2　an acknowledgement transaction. If the transaction does not have a correct sequence

3　number, the transaction is dropped, since an incorrect sequence number means that earlier

4　transactions have been dropped in the system. If the destination node cannot generate a

5　response (or the transaction is a response transaction) the destination node simply sends

6　an acknowledgement transaction to the source node. In either case, the destination node

7　resets the flag bit in the receive_seqid table indicating that the acknowledgement (or

8　response) has been sent. The destination node sends acknowledgement transactions for

9　transactions received from a particular node, path and flow control class, in order.

10　　　If a source node does not receive an acknowledgement transaction within a

11　predetermined time, the source node sends a probe request transaction along an alternate

12　path (preferably an alternate physical path). The probe request transaction contains the

13　source node identification, the path number, the flow control class, the sequence number

14　of the timed-out transaction, and the sequence number of the last transaction that is

15　pending. The destination node takes the information contained in the probe request

16　transaction and determines if the destination node has already responded to the timed-out

17　transaction. If the destination node has already responded to the timed-out transaction,

18　the destination node indicates so in a probe request response along with the sequence

19　number of the last transaction that the destination node has received. This probe request

20　response is sent along an alternate path. The probe request transaction, as well as the

21　corresponding probe request response, may then be used for acknowledgement purposes.

22　When the source node receives acknowledgement to the probe request transaction, the

23　source node resumes retransmission starting with the transaction after the last sequence

24　number received by the destination node, if any. From this point on, neither the source

25　node nor the destination node use the path where the problem occurred to receive a

26　transaction or to send out an acknowledgement.

27　**Description of the Drawings**

28　　　The detailed description will refer to the following figures, in which like numbers

29　refer to like elements, and in which:

30　　　Figure 1 is a diagram of a multiprocessor computer system that employs a static

31　end to end retransmit apparatus and method;

32　　　Figure 2 is a further block diagram of a system of Figure 1;

33　　　Figure 3 illustrates a sequence identification table used with the apparatus of

34　Figure 1; and

1    Figures 4-8 are flowcharts of operations of the apparatus of Figure 1.

2    **Detailed Description**

3    A static end to end retransmit protocol is implemented by an end to end retransmit

4    apparatus and method employed in a multiprocessor computer system. Figure 1 is a

5    block diagram of a multiprocessor computer system 10 that employs such an apparatus.

6    In Figure 1, a source node 11 is coupled to a destination node 12 through alternate paths

7    20. The source node 11 is also coupled to a retransmit buffer 13 and a send_seqid table

8    15 and a receive_seqid table 16. The destination node 12 is coupled to a receive buffer 17

9    and a receive_seqid table 18 and a send_seqid table 19. The designation of the nodes 11

10   and 12 is arbitrary, and for means of illustration. In the system 10, both the nodes 11 and

11   12 may send and receive transactions and hence both the nodes 11 and 12 may be any

12   source or destination nodes. The nodes 11 and 12 may be any nodes in the

13   multiprocessor computer system 10, such as CPU or memory nodes or I/O hub chips.

14   The paths 20 may be distinct physical paths or virtual paths. The source node 11 sends

15   transactions to the destination node 12 along one of the physical paths 20 and receives

16   responses or acknowledgements from the destination node 12 along one of the physical

17   paths 20. Transmissions sent from the source node 11 to the destination node 12 may be

18   temporarily placed in the retransmit buffer 13. Similarly, responses and

19   acknowledgements from the destination node 12 to the source node 11 may be

20   temporarily placed in the receive buffer 17. The send_seqid tables 15 and 19 and the

21   receive_seqid tables 16 and 18 may be used to store information related to the

22   transactions such as the sequence number (or sequence ID) of each transaction, response

23   or acknowledgement.

24   Figure 2 is a block diagram of the microprocessor computer system 10 of Figure 1

25   showing additional details of operation of the end to end retransmit apparatus. The nodes

26   11 and 12 are connected through a number of cross-bar type routing chips. In the

27   illustrated example, the cross-bar chips 34, 36, and 38 are used to connect the nodes 11

28   and 12. However, fewer or more cross-bar chips could be used.

29   The nodes 11 and 12 are connected by two paths, P1 and P2. The path P1 is

30   designated by links 21 - 24. The path P2 is designated by links 25 - 28. Any of the links

31   21 – 28 in the paths P1 and P2 may fail. For example, the link 28 (in path P2 from the

32   source node 11 to the destination node 12) may fail. Thereafter, any transaction the

33   source node 11 sends (or has sent) to the destination node 12 over the link 28 and the path

34   P2 not arrive at the destination node 12, and hence will not be acknowledged by the

1    destination node 12.  The source node 11 will eventually time out on the oldest non-

2    acknowledged transaction.  The source node 11 will thenceforth stop using the path P2

3    for sending any subsequent normal transactions.  In particular, the source node 11 will

4    deconfigure the path P2 and will stop accepting any acknowledgements that are sent to

5    the source node 11 over the path P2.  However, the source node 11 may continue to

6    receive normal transactions over the path P2.  The source node 11 may also send a probe

7    request along the path P1, for example, over the links 21 and 23.  The destination node 12

8    may respond, using the path P1, with the sequence number of the last transaction received

9    by the destination node 12 from the source node 11 over the path P2.  The destination

10   node 12 then stops receiving any normal transactions along the path P2.  The

11   deconfigured path P2 may be indicated by use of a separate status bit, or by simply

12   reserving a sequence number (for example, all 0's).

13          The source node 11 may attempt to determine if the failed path P2 is still open.

14   For example, the unacknowledged transaction may have been the result of a transient

15   error, in which case the path P2 may still be available for sending and receiving

16   transactions, including acknowledgements.  After receiving the response to the probe

17   request, the source node 11 may send a plunge request along the failed path P2 and flow

18   control class to the destination node 12.  The plunge request indicates the sequence

19   number of the first transaction the source node 11 will retransmit if the path P2 is re-

20   established.  On receiving the plunge request, the destination node 12 may re-establish the

21   path P2.  The destination node 12 then initiates a response for the plunge request.  Since

22   the plunge request itself may be in the response flow control class, the destination node

23   12 may use a flag bit in the receive_seqid table 18 to send the plunge request response

24   when space exists in the receive buffer 17.  Once the source node 11 receives the

25   response to the plunge request, the source node 11 can start using the path P2 for normal

26   transactions.  If the source node 11 does not receive a response to the plunge request, the

27   source node 11 does not use the path P2 until maintenance software guarantees that the

28   path P2 has been re-established.  In an embodiment, the source node 11 may retry the

29   determination of the existence of the path P2 by periodically sending plunge requests to

30   the destination node 12.

31          In the multiprocessor computer system 10 shown in Figures 1 and 2, each of the

32   nodes 11 and 12 keeps track of transactions the node has sent over time to the other node,

33   as well as every transaction the node has received from the other node, along each active

34   path for each flow control class.  To accomplish this tracking function, two data structures

1    exist.   A send_seqid, represents the sequence number (sequence ID) for the last

2    transaction sent by the source node 11 to the destination node 12 along any given active

3    path and for each flow control class.  A receive_seqid represents the sequence number of

4    the last transaction that the destination node 12 received and for which the destination

5    node 12 sent an acknowledgement (ACK) back to the source node 11, for each node,

6    along the active path, and for each flow control class.  Each node (destination node 12 for

7    send_seqid or source node 11 for receive_seqid) can operate over multiple paths.  All

8    nodes in one full control class may operate over the same number of paths.  For example,

9    the system 10 may have four alternate active paths between any two CPU/memory nodes,

10   but one active path to or from an I/O hub chip.  The system 10 does not require distinct

11   physical paths between any source-destination nodes.  For example, the system 10 may

12   comprise four active paths with two active paths sharing a physical path.

13         The flow control class refers to the type of transactions being sent over the

14   paths/links in the system 10.  For example, a memory read request may be in flow control

15   class A and a write request in flow control class B.  A path that is available to one flow

16   control class may not be available to different flow control class.

17         Every transaction that is sent from the source node 11 to the destination node 12 is

18   also put into the retransmit buffer 13.  When the transaction gets an acknowledgement

19   from the destination node 12, the transaction is removed from the retransmit buffer 13.

20   The acknowledgement can be piggy-backed with an incoming transaction and/or a special

21   transaction.  No acknowledgement is necessary for an acknowledgement transaction.  If a

22   transaction does not get an acknowledgement within a predetermined time, recovery

23   actions may be taken.  The destination node 12 may wait to gather several transactions for

24   a given source node 11 before generating an explicit acknowledgement transaction, while

25   trying to ensure that such a delay will not generate any recovery actions at the source

26   node 11.  This delay helps conserve bandwidth by avoiding explicit acknowledgement

27   transactions as much as possible.

28         When the source node 11 sends a transaction to a destination node 12, the source

29   node 11 gets the sequence number from the send_seqid table 15, checks that the

30   transaction with the sequence number is pending to the same destination node 12, and

31   sends the transaction to the destination node 12 while placing the transaction in the

32   retransmit buffer 13.  The source node 11 then increments the corresponding sequence

33   number in the send_seqid table 15.  When the destination node 12 receives the

34   transaction, the destination node 12 queues the transaction in the receive buffer 17.  If the

1 transaction is of a request type, and the destination node 12 can generate a response
2 within the time out period, the destination node 12 sends a response to the source node
3 11, which acts as an implicit acknowledgement. The destination node 12 then checks the
4 receive_seqid table 18 to see if the transaction the destination node 12 received has the
5 correct sequence number. If the transaction has the correct sequence number, the
6 destination node 12 updates the corresponding entry in the receive_seqid table 18, and
7 sets a flag bit indicating that the destination node 12 needs to send an acknowledgement
8 transaction. If the transaction does not have a correct sequence number, the transaction is
9 dropped, since an incorrect sequence number means that earlier transactions have been
10 dropped in the system 10. If the destination node 12 cannot generate a response (or the
11 transaction is a response transaction) the destination node 12 simply sends an
12 acknowledgement transaction to the source node 11. In either case, the destination node
13 12 resets the flag bit in the receive_seqid table 18 indicating that the acknowledgement
14 (or response) has been sent. The destination node 12 sends the acknowledgment for all
15 transactions received from the source node 11, along a given active path, and flow control
16 class, in order.

17 If the source node 11 does not receive an acknowledgement transaction within a
18 predetermined time, the source node 11 sends a probe request transaction along an
19 alternate path (preferably an alternate physical path). The probe request transaction
20 contains the source node identification, the path number, the flow control class, and the
21 sequence ID of the timed-out transaction, and the sequence number of the last transaction
22 that is pending in the retransmit buffer 13. The destination node 12 takes the information
23 contained in the probe request transaction and determines if the destination node 12 has
24 already responded to the timed-out transaction. If the destination node 12 has already
25 responded to the timed-out transaction, the destination node 12 indicates so in a probe
26 request response along with the sequence number of the last transaction which the
27 destination node 12 has received. The probe request response is sent along an alternate
28 path. The probe request transaction, as well as the corresponding probe request response,
29 may then be used for acknowledgement purposes. When the source node receives
30 acknowledgement to the probe request transaction, the source node resumes
31 retransmission starting with the transaction after the last sequence number received by the
32 destination node 12, if any. From this point on, neither the source node 12 nor the
33 destination node 13 use the path where the problem occurred to receive a transaction or to
34 send out an acknowledgement.

1    Figure 3 illustrates an exemplary receive_seqid table 18 used with the system 10

2    of Figures 1 and 2.  The table 18 comprises entries for each pending transaction.  Each

3    entry in the table comprises a sequence number (sequence ID) and a flag bit that is set to

4    either one or zero.  The sequence number includes the node identification of the

5    transaction source node, the path, and the flow control class.  The sequence number may

6    be set to all zeros to indicate that a particular path is deconfigured.  The flag bit may be

7    set to zero to indicate that an acknowledgement has already been sent, or may be set to

8    one to indicate that an acknowledgement still needs to be sent.  The flag bit set to one

9    may also indicate, in the case of a deconfigured path, the need to send a plunge request

10   response.  The send_seqid table 15 may be similar to the receive_seqid table 18

11   illustrated.  Each mode in the computer system 10 may include a receive_seqid and a

12   send_seqid table.  Alternatively, separate data structures similar to that illustrated may be

13   provided for each flow control class.

14       Figures 4 – 8 illustrate various operations of the system 10 shown in Figures 1 and

15   2.  The operations described below assume that transactions are sent along one of the

16   paths P1 or P2 from the source node 11 to the destination node 12, and that responses and

17   acknowledgements are returned.  The transactions may be in one or more flow control

18   classes, and flow control class F will be used by way of example.

19       Figure 4 shows the steps involved in an operation 100 for sending a regular, or

20   normal, transaction in flow control class F, along path P2 from the source node 11 to the

21   destination node 12.  The operation 100 begins in block 105.  In block 110, the

22   destination node 11 retrieves a sequence number for a pending (designated) transaction

23   from the send_seqid table 15.  In block 115, the node 11 determines if the retransmit

24   buffer 13 contains a transaction with the same sequence number as that retrieved in block

25   110, for a transaction to the destination node 12, using path P2, and flow control class F.

26   Should the retransmit buffer 13 contain such a transaction, the operation moves to block

27   120, and the operation 100 waits for a time to allow the transaction in the retransmit

28   buffer 13 to be cleared.  Such a wait is necessary to avoid assigning a same sequence

29   number to more than one transaction.  Alternatively, the source node 11 may send the

30   transaction to the destination node 12 along another available path, such as the path P1.

31   Following an appropriate wait, or if another path is available, the operation 100 returns to

32   block 115 to determine if the retransmit buffer 13 contains a pending transaction with the

33   same sequence number, path and flow control class as the transaction designated in block

34   110.

1    In block 115, if the retransmit buffer 13 does not contain a transaction with the

2    same sequence number, path and flow control class as the designated transaction, the

3    operation 100 moves to block 125, and the retrieved sequence number is attached to the

4    designated transaction. The designated transaction is then sent along the path P2 to the

5    destination node 12. In block 130, the thus-sent designated transaction is placed in the

6    retransmit buffer 13. In block 135, the operation 100 ends.

7    Figure 5 illustrates an operation 200 that occurs when a transaction in the

8    retransmit buffer 13 times out. The operation 200 begins in block 205, with a transaction

9    that has timed out. In block 210, the source node 11 retrieves the sequence number of the

10   last (most recent) transaction sent to the destination node 12 along the path P2 and the

11   flow control class F as the transaction that timed out. For example, a transaction with a

12   sequence number (ID) of 2 may have timed out, and the most recent transaction sent to

13   the destination node 12 has a sequence number of 5 (both transactions 2 and 5 in the same

14   flow control class, and sent along the same path). The source node 11 retrieves the

15   sequence number 5. In block 215, the source node 11 sends a probe request transaction to

16   the destination node 12 along alternate path P1 with the sequence number of the timed-

17   out transaction and the sequence number of the last transaction sent, for the flow control

18   class F. The source node 11 then deconfigures the path P2, block 220, for both sending

19   transactions and receiving acknowledgments, updating the sequence number to 0. This

20   step prevents using a path over which a timed-out transaction was to have been

21   transmitted.

22   Once the probe request has been sent and the path deconfigured, the source node

23   11 waits to receive a probe response from the destination node 12. The wait can be

24   established to account for a two-way transmission, plus some additional time in case of

25   delays and tolerance errors. For example, the wait may be four times the design two-way

26   transmission time. In block 225, the source node 11 determines if the probe response has

27   been received. If the probe response has not been received, the source node 11 may

28   determine if the probe response has timed out, block 230. If the probe response has not

29   timed out, the operation 200 returns to block 225, and the source node 11 continues to

30   wait for the probe response. If the probe response has timed out, the operation 200 moves

31   to block 235, and the source node 11 determines if another path (for example a path P3) is

32   available for sending the probe response to the destination node 12. If another path is

33   available, the operation 200 returns to block 215. If another path is not available, then an

1    error condition exists and the operation moves to block 270 and ends, indicating to

2    management software that possible uncorrectable errors exist.

3         In block 225, if the source node 11 determines that the probe response has been

4    received from the destination node 12, the operation 200 moves to block 240, and the

5    source node 11 resumes transmission of all transactions in the retransmit buffer 13 for the

6    destination node 12, along path P2 and flow control class F, for which the destination

7    node 12 has not provided an acknowledgement.

8         At this point, the operation 200 may end as shown by the dashed line connecting

9    blocks 240 and 270. Alternatively, the operation 200 may continue to block 245, and the

10   source node 11 sends a plunge transaction to the destination node 12 along the path P2 in

11   the flow control class F (that is, the plunge transaction follows the transaction

12   corresponding to the last sequence number sent) indicating the sequence number where

13   the source node 11 will begin transmission should the failed path P2 be re-established.

14   The source node 11 next updates the sequence number in the send_seqid table 15. The

15   source node 11 then waits for a response to the plunge transaction. Since the response to

16   the plunge transaction may be in the same flow control class as the plunge transaction

17   itself, the destination node 12 may simply use the flag bit in the receive_seqid table 18 to

18   indicate when to send the plunge response if space does not exist in the receive buffer 17.

19        In block 250, the source node 11 determines if the plunge response has been

20   received. If the plunge response has not been received, the source node 11 determines if

21   the plunge transaction has timed out, block 255. If the plunge transaction has not timed

22   out, the source node 11 continues to wait for the plunge response. If the plunge

23   transaction has timed out, the source node 11 waits for either a software intervention to

24   indicate that the path P2 is available, or may optionally wait for a longer time before

25   sending another plunge response along the path P2, block 260. In block 250, if the source

26   node 11 determines that the plunge response has been received, the operation 200 moves

27   to block 265, and the source node 11 reconfigures the path P2 for normal transactions.

28   The operation then returns to block 240, and any pending transactions in the retransmit

29   buffer 13 for which an acknowledgement has not been received are sent to the destination

30   node 12.

31        Figure 6 is a flow chart of an operation 300 for at the destination node 12 for

32   receiving and responding to a plunge transaction. The operation 300 begins in block 305.

33   In block 310 the destination node 12 receives a plunge transaction. In block 315, the

34   destination node 12 determines if the path and flow control class from the source node 11

1 is deconfigured (turned off). If the path and flow control class is deconfigured, the
2 operation 300 moves to block 320 and the destination node 12 sends a plunge response
3 back to the source node 11 along the same path the plunge transaction used. In block
4 320, if the destination node 12 cannot send a plunge response immediately (because, for
5 example, the receive buffer 17 is full), the destination node 12 may use the flag bit in the
6 receive_seqid table 18 to indicate when to send the plunge response to the source node
7 11. The destination node 12 then reconfigures the path P2 for sending and receiving
8 normal transactions. The operation 300 then ends, block 330.

9 In block 315, if the path is not deconfigured, the operation 300 moves to block
10 325, and the destination node 12 notes an uncorrectable error, notifies the source node 11,
11 and deconfigures the path P2. The operation 300 then ends, block 330.

12 Figure 7 is a flow chart illustrating an operation 400 at the destination node 12 for
13 receiving and responding to a probe transaction. The operation 400 begins in block 405.
14 In block 410, the destination node 12 receives a probe transaction from the source node
15 11. In block 415, the destination node 12 looks up the latest entry from the source node
16 11, for the path P2 and flow control class F, in the receive_seqid table 18 or the receive
17 buffer 17. The operation 400 then moves to block 420, and the destination node 12
18 determines if either of the following two conditions is met: (1) the latest sequence number
19 corresponds to the time-out sequence number, minus one; or (2) the latest sequence
20 number lies between the timed-out sequence number and the last sequence number sent
21 from the source node 11. If neither condition is met, the operation 400 moves to block
22 430, and an uncorrectable error condition is declared.

23 In block 420, if the destination node 12 determines that either condition is met, the
24 operation 400 moves to block 425, and the destination node 12 sends a probe response
25 along the alternate path (e.g., path P1) over which the probe transaction was transmitted.
26 The destination node 12 also indicates if the timed-out transaction was received, and
27 sends the source node 11 the sequence number of the last transaction acknowledged.
28 Following blocks 425 and 430, the operation 400 ends, block 435.

29 Figure 8 is a flow chart showing a normal transaction operation 500 at the
30 destination node 12. In the example illustrated, the source node 11 sends a normal
31 transaction to the destination node 12 along the path P2 in flow control class F. The
32 operation 500 begins in block 505. In block 510, the destination node 12 receives a
33 normal (regular) transaction. In block 515, the destination node 12 looks up in the
34 receive_seqid table 18/the receive buffer 17 the sequence number, path and flow control

1    class for the last transaction received from the source node 11.   In block 520, the

2    destination node 12 determines if the sequence number for the newly received normal

3    transaction is subsequent to the sequence number of the last transaction as determined in

4    block 515.  If the sequence number of the new transaction is subsequent to that of the last

5    transaction, the operation 500 moves to block 525 and the destination node 12 determines

6    if the path P2 is still configured.   If the path P2 is still configured, the operation 500

7    moves to block 530, and the destination node 12 accepts the normal transaction.   The

8    destination node 12 then updates the acknowledgement to be sent, updates the

9    receive_seqid table 18, and send the acknowledgement to the source node 11 along the

10   path P2.  If in either block 520 or 525, the condition is not met, the operation moves to

11   block 535, and the transaction is dropped.  The operation 500 then moves to block 540

12   and ends.